

**Ninth Conference
of the
European Chapter
of the
Association for
Computational Linguistics**

8–12 June 1999
University of Bergen
Bergen, Norway

Published by the Association for Computational Linguistics

The conference was sponsored by

LINGSOFT

University of Bergen (Humanities Faculty)

Bergen University Fund

Norwegian Ministry of Education, Research and Church affairs

Nordisk Språkteknologi AS (Nordic Language Technology)

Ministry of Culture

©Copyright 1999, Association for Computational Linguistics

Order copies of this and other ACL publications from:

Morgan Kaufmann Publishers

340 Pine Street, 6th floor

San Francisco, CA 94104, USA

+1 800 745 7323

Preface

This volume contains the papers prepared for the Ninth Conference of the European Chapter of the Association for Computational Linguistics, held 8–12 June in Bergen, Norway.

The conference programme features invited talks, tutorials, submitted papers in both main and student sessions and a poster and demo session.

The main session Programme Committee received 99 abstracts from 20 countries around the world. Every paper was reviewed by at least three reviewers via a web-based interface which preserved anonymity. The fairness of the process is we think demonstrated by the broad spread of institutions and countries represented in the papers which were accepted and are printed here: 29 papers from 12 countries, with no country having more than 6 papers. The topics of the accepted papers cover a wide range of topics, and taken together we think they present an exciting and up-to-date sample of the best work in our field at the present time.

Our Area Chairs worked extremely hard in recruiting the Programme Committee, allocating papers to reviewers and encouraging the completion of reviews on time, and then joined the Programme Chairs for an intensive all-day meeting in Edinburgh where the final decisions were made. We are extremely grateful to both the Area Chairs and the Programme Committee for their hard work.

We used a new web-based approach to managing the paper submission and reviewing process this year, with authors registering an intention to submit electronically, and web-based reviewing processes, while at the same time retaining paper submission of the manuscripts themselves. We've had a lot of positive feedback about the reviewing system from the Programme Committee. We hope the system worked for authors as well.

Pulling together a meeting on this scale is a major effort, and thanks are due to many people for making it happen: Bruce Croft and Wolfgang Wahlster for the invited talks; Walter Daelemans for organising the tutorials and Robert Dale, Ronen Feldman, Adam Kilgariff, Adwait Ratnaparkhi, Ehud Reiter and Michael Rundell for preparing and delivering them; the EACL President (Donia Scott) and past President (John Nerbonne) for their helpful suggestions; Gertjan van Noord for organising the post-conference workshops and Jonas Kuhn and Avro Voutilainen for organising the student sessions.

We are also very grateful to the Local Organisation Committee Chair, Koenraad de Smedt, for his patient cooperation with us over many issues, large and small, as well as his stirring efforts in making sure we would actually *have* a conference at which these papers could be presented.

We also are indebted to Robert Inder made his web-based registration and reviewing system available to us, and Richard Tobin helped us adapt it to our needs.

Local help in organising the Programme Committee meeting; with incoming and outgoing mail and with producing these Proceedings was provided by Janet Forbes and Margaret McMillan, to whom many thanks.

Henry S. Thompson, Alex Lascarides
Programme Committee Chairs

Programme Committee (Main Session)

Programme Chair: Programme Co-Chair:

Henry S. Thompson (University of Edinburgh) Alex Lascarides (University of Edinburgh)

Area Chairs:

John Carroll (University of Sussex) Rob Gaizauskas (Sheffield University)
Jan Hajic (Charles University) Mats Rooth (University of Stuttgart)
Paul Taylor (University of Edinburgh) Lyn Walker (AT&T Bell Labs)

Programme Committee Members:

Ann Abeille (University of Paris 7) Steve Abney (AT&T Bell Labs)
Elisabeth Andre (DFKI, Saarbrücken) Susan Armstrong (ISSCO, Univ. of Geneva)
Doug Arnold (University of Essex) Paolo Baggia (CSELT)
Ken Beesley (Xerox Research Centre Europe) Bill Black (UMIST)
Gosse Bouma (University of Groningen) Eric Brill (CLSP Johns Hopkins Univ.)
Lynne Cahill (University of Sussex) Nicoletta Calzolari (CNR, Pisa)
Bob Carpenter (Lucent Technologies) Kenneth W. Church (AT&T Bell Labs)
Dick Crouch (University of Nottingham) Fabio Ciravegna (IRST)
Michael Collins (University of Pennsylvania) Ann Copestake (Stanford University)
Morena Danieli (CSELT) Judy Delin (University of Stirling)
Ted E. Dunning (Aptex) Jan van Eijck (CWI)
Michael Elhadad (Ben Gurion University) Tomaz Erjavec (Ljubljana Univ.)
Roger Evans (University of Brighton) Frank van Eynde (Katholieke Universiteit Leuven)
Dan Flickinger (Stanford University) Claire Gardent (University of Saarbrücken)
Gregory Greffentette (Xerox Research Centre Europe) Ralph Grishman (New York University)
Udo Hahn (University of Freiburg) Jiri Hana (Charles Univ., Prague)
Djoerd Hiemstra (University of Twente) Julia Hirschberg (AT&T Bell Labs)
Barbora Hladka (Charles Univ., Prague) Beryl Hoffman (Marymount University)
David Israel (SRI International) Karen Jensen (Microsoft)
Arne Johnson (University of Linköping) Mark Johnson (Brown University)
Douglas A. Jones (Dept. of Defense) Ewan Klein (University of Edinburgh)
Marcus Kracht (Free University of Berlin) Hans-Ulrich Krieger (DFKI, Saarbrücken)
Geert-Jan Kruijff (Charles University, Prague) Ivana Kruijff (Charles University, Prague)
Gina Levow (University of Maryland) Hang Li (C&C Media Research Labs NEC Corp.)
Marc Light (MITRE Corporation) Suresh Manandhar (University of York)
Inderjeet Mani (MITRE Corporation) Mitch Marcus (Univ. of Pennsylvania)
Andrei Mikheev (Harlequin) David Milward (SRI Cambridge)
Marc Moens (University of Edinburgh) Johanna Moore (University of Edinburgh)
Christine Nakatani (Lucent Technologies) Mark-Jan Nederhof (DFKI, Saarbrücken)
Guenter Neumann (DFKI, Saarbrücken) Anton Nijholt (University of Twente)
Gertjan van Noord (University of Groningen) Kemal Oflazer (New Mexico State University)
Karel Oliva (Univ. of Saarbrücken) Chris Paice (University of Lancaster)
Richard Power (University of Brighton) Victor Poznanski (Sharp Laboratories of Europe)
Steve Pulman (Cambridge University) Owen Rambow (CogenTex Inc)
Ehud Reiter (University of Aberdeen) Norbert Reithinger (DFKI, Saarbrücken)
Eileen Riloff (University of Utah) Craig Roberts (Ohio State University)
Bill Rounds (University of Michigan) Donia Scott (University of Brighton)
Harold Somers (UMIST) Richard Sproat (AT&T Bell Labs)
Tomasz Strzalkowski (General Electric) Marc Swerts (Institute for Perception Research)
John Tait (University of Sunderland) Louis des Tombe (Utrecht University)
Marc Villain (MITRE Corporation) Juergen Wedekind (University of Stuttgart)
David Weir (University of Sussex) Janyce Wiebe (New Mexico State University)
Dekai Wu (Hong Kong Univ. of Science and Tech.) Antonio Zampolli (CNR, Pisa)
Dan Zeman (Charles University, Prague) Michael Zock (LIMSI)
Ingrid Zukerman (Monash University)

Student Session Preface

The student sessions have become an integral part of the EACL conferences. They provide an invaluable opportunity for young researchers to present their work to the community and receive feedback for future activities. In this spirit, we kept the tradition of encouraging students to submit not only papers presenting completed work (like for the main sessions), but also reports on work in progress.

We received 17 submissions from 8 countries, of which we accepted 8 papers for presentation and two as reserve papers. Each submission was reviewed at least by two student reviewers and one faculty reviewer.

We would like to thank the 31 reviewers for their fair assessment and their detailed comments, which we think were of great help for the student authors. Their names and affiliations are listed below.

We want to thank the Student Sessions Co-chairs of this year's ACL conference Melanie Baljko (Toronto, Canada) and Anna Korhonen (Cambridge, UK), and the Student Session Chairs of previous ACL and COLING-ACL conferences, in particular Pamela W. Jordan (University of Pittsburgh, USA), Maria Milosavljevic (CSRIRO, Australia), and Dragomir R. Radev (Columbia University, USA), for providing supporting material and for their cooperation. Finally, we want to thank the organizers of the main conference, more specifically Koenraad de Smedt (Bergen, Norway), Alex Lascarides (Edinburgh, UK), Giorgio Satta (Padova, Italy), and Henry Thompson (Edinburgh, UK), and the members of the European Chapter of ACL for their cooperation and help.

Jonas Kuhn, Atro Voutilainen

Programme Committee (Student Session)

Programme Chair: Programme Co-Chair:

Jonas Kuhn (Stuttgart, Germany) Atro Voutilainen (Helsinki, Finland)

Programme Committee Members:

Tilman Becker (Saarbrücken, Germany)	Gosse Bouma (Groningen, Netherlands)
Stefanie Dipper (Stuttgart, Germany)	Björn Gambäck (SICS, Sweden)
Phil Harrison (London, UK)	Joris Hulstijn (Utrecht, Netherlands)
Timo Järvinen (Helsinki, Finland)	Frank Keller (Edinburgh, UK)
Alexandra Kinyon (Paris, France/Pennsylvania, USA)	Jacques Koreman (Saarbrücken, Germany)
Peter Krause (Stuttgart, Germany)	Jan van Kuppevelt (Stuttgart, Germany)
Kordula de Kuthy (Saarbrücken, Germany)	Maria Lapata (Edinburgh, UK)
Scott McDonald (Edinburgh, UK)	Guido Minnen (Sussex, UK)
Bernd Möbius (Stuttgart, Germany)	Christof Monz (Amsterdam, Netherlands)
Mark-Jan Nederhof (Saarbrücken, Germany)	Nicolas Nicolov (Sussex, UK)
Lluís Padró (Barcelona, Spain)	Jussi Piitulainen (Helsinki, Finland)
Peter Poller (Saarbrücken, Germany)	Detlef Prescher (Stuttgart, Germany)
Stefan Riezler (Stuttgart, Germany)	Antje Roßdeutscher (Stuttgart, Germany)
Henk Schotel (Maastricht, Netherlands)	Mariet Theune (Eindhoven, Netherlands)
Gökhan Tür (SRI, Menlo Park, CA, USA)	Jennifer Venditti (AT&T Bell Labs, USA)
Anssi Yli-Jyrä (Helsinki, Finland)	

Programme Committee (Poster Session)

Programme Chair:

Giorgio Satta (University of Padua)

Programme Committee Members:

Breck Baldwin (University of Pennsylvania)	Tilman Becker (DFKI, Saarbrücken)
Mari Broman Olsen (University of Maryland)	Federica Busa (Brandeis University)
Marie-Helene Candito (Universite Paris 7)	Thierry Declerck (DFKI, Saarbrücken)
Luca Dini (CELI, SNS, Italy)	Christy Doran (University of Brighton)
Jason Eisner (University of Pennsylvania)	Karin Harbusch (University of Koblenz)
Mark Hepple (University of Sheffield)	Elena Not (IRST, Italy)
Carolyn Penstein Rosé	Fabio Pianesi (IRST, Italy)
Emanuele Pianta (IRST, Italy)	Owen Rambow (CogenTex Inc.)
Norbert Reithinger (DFKI, Saarbrücken)	Srinivas Bangalore (AT&T Bell Labs)
Pasi Tapanainen (University of Helsinki)	Jakub Zavrel (Tilberg University)

Table of Contents

Andrei Mikheev, Marc Moens and Claire Grover <i>Named Entity Recognition without Gazetteers</i>	1
Richard Power <i>Generating Referring Expressions with a Unification Grammar</i>	9
Didier Bourigault and Christian Jacquemin <i>Term Extraction and Term Clustering: An Integrated Platform for Computer-Aided Terminology</i>	15
Ian M. O'Neill and Michael F. McTear <i>An Object-Oriented Approach to the Design of Dialogue Management Functionality</i>	23
Maria Lapata, Scott McDonald and Frank Keller <i>Determinants of Adjective-Noun Plausibility</i>	30
Miriam Eckert and Michael Strube <i>Resolving Discourse Deictic Anaphora in Dialogues</i>	37
Suzanne Stevenson and Paola Merlo <i>Automatic Verb Classification Using Distributions of Grammatical Features</i>	45
Pierre Boullier <i>Chinese Numbers, MIX, Scrambling, and Range Concatenation Grammars</i>	53
Glyn Morrill <i>Geometry of Lexico-Syntactic Interaction</i>	61
Franz Josef Och <i>An Efficient Method for Determining Bilingual Word Classes</i>	71
Inderjeet Mani, Therese Firmin, David House, Gary Klein Beth Sundheim and Lynette Hirschman <i>The TIPSTER SUMMAC Text Summarization Evaluation</i>	77
Tim Fernando <i>Ambiguous Propositions Typed</i>	86
Rila Mandala, Takenobu Tokunaga and Hozumi Tanaka <i>Complementing WordNet with Roget's and Corpus-based Thesauri for Information Retrieval</i>	94
Fabio Ciravegna and Alberto Lavelli <i>Full Text Parsing using Cascades of Rules: An Information Extraction Perspective</i>	102
Simone Teufel, Jean Carletta and Marc Moens <i>An Annotation Scheme for Discourse-level Argumentation in Research Articles</i>	110
Thorsten Brants <i>Cascaded Markov Models</i>	118
Dale Gerdemann and Gertjan van Noord <i>Transducers from Rewrite Rules with Backreferences</i>	126
Giorgos S. Orphanos and Dimitris N. Christodoulakis <i>POS Disambiguation and Unknown Word Guessing with Decision Trees</i>	134

Maria Wolters and Mathias Kirsten <i>Exploring the Use of Linguistic Features in Domain and Genre Classification</i>	142
Miguel A. Alonso, David Cabrero, Eric de la Clergerie and Manuel Vilares <i>Tabular Algorithms for TAG Parsing</i>	150
Efstathios Stamatatos, Nikos Fakotakis and George Kokkinakis <i>Automatic Authorship Attribution</i>	158
Guido Minnen <i>Selective Magic HPSG Parsing</i>	165
Erik F. Tjong Kim Sang and Jorn Veenstra <i>Representing Text Chunks</i>	173
Hiroyuki Shinnou <i>Detection of Japanese Homophone Errors by a Decision List Including a Written Word as a Default Evidence</i>	180
John Chen, Srinivas Bangalore and K. Vijay-Shanker <i>New Models for Improving Supertag Disambiguation</i>	188
Kiyotaka Uchimoto, Satoshi Sekine and Hitoshi Isahara <i>Japanese Dependency Structure Analysis based on Maximum Entropy Models</i>	196
Atro Voutilainen <i>An Experiment on the Upper Bound of Interjudge Agreement: The Case of Tagging</i>	204
Fumiyo Fukumoto and Yoshimi Suzuki <i>Word Sense Disambiguation in Untagged Text Based on Term Weight Learning</i>	209
John Carroll, Nicolas Nicolov, Olga Shaumyan, Martine Smets and David Weir <i>Parsing with an Extended Domain of Locality</i>	217

Student Session

Gabriela Cavaglià <i>The Development of Lexical Resources for Information Extraction from Text Combining WordNet and Dewey Decimal Classification</i>	225
Donna K. Byron and Joel R. Tetreault <i>A Flexible Architecture for Reference Resolution</i>	229
Patrick Caudal <i>Result States and the Lexicon: The Proper Treatment of Event Structure</i>	233
Daniel S. Paiva <i>Investigating NLG Architectures: Taking Style into Consideration</i>	237
Justin Picard <i>Finding Content-bearing Terms Using Term Similarities</i>	241
Dimitrios Kokkinakis and Sofie Johansson Kokkinakis <i>A Cascaded Finite-State Parser for Syntactic Analysis of Swedish</i>	245
Patrice Lopez <i>Repair Strategies for Lexicalized Tree Grammars</i>	249
Veit Reuer <i>Dialogue Processing in a CALL-System</i>	253
Yan Zuo <i>Focusing on Focus: A Formalization</i>	257
Aline Villavicencio <i>Representing a System of Lexical Types Using Default Unification</i>	261

Poster Session

Izaskun Aldezabal, Inaki Alegria, Olatz Ansa, Jose Mari Arriola, Nerea Ezeiza, Itziar Aduriz and Alexander Da Costa <i>Designing Spelling Correctors for Inflected Languages Using Lexical Transducers</i>	265
Hans Argenton and Anke Feldhaus <i>The Treegram Index—An Efficient Technique for Retrieval in Linguistic Treebanks</i>	267
John Carroll, Guido Minnen, Darren Pearce, Yvonne Canning, Siobhan Devlin and John Tait <i>Simplifying Text for Language-Impaired Readers</i>	269
Nigel Collier, Hyun Seok Park, Norihiro Ogata, Yuka Tateishi, Chikashi Nobata, Tomoko Ohta, Tateshi Sekimizu, Hisao Imai, Katsutoshi Ibushi and Jun-ichi Tsujii <i>The GENIA Project: Corpus-based Knowledge Acquisition and Information Extraction from Genome Research Papers</i>	271
Crit Cremers <i>A Note on Categorical Grammar, Disharmony and Permutation</i>	273
Johann Gamper <i>Encoding a Parallel Corpus for Automatic Terminology Extraction</i>	275
Adam Kilgarriff <i>95% Replicability for Manual Word Sense Tagging</i>	277
Torbjörn Lager <i>μ-TBL Lite: A Small, Extendible Transformation-Based Learner</i>	279
John Nerbonne, Wilbert Heeringa and Peter Kleiweg <i>Comparison and Classification of Dialects</i>	281
Frank Schilder <i>Pointing to Events</i>	283
Mark Stevenson <i>A Corpus-Base Approach to Deriving Lexical Mappings</i>	285
José Relano Gil, Daniel Tapias, Maria C. Gancedo, Marcela Charfuelan and Luis A. Hernández <i>Robust and Flexible Mixed-Initiative Dialogue for Telephone Services</i>	287
Kaili Müürisep <i>Determination of Syntactic Functions in Estonian Constraint Grammar</i>	291

Author Index

Miguel A. Alonso	150	Michael F. McTear	23
Srinivas Bangalore	188	Paola Merlo	45
Pierre Boullier	53	Andrei Mikheev	1
Didier Bourigault	15	Guido Minnen	165
Thorsten Brants	118	Marc Moens	1, 110
David Cabrero	150	Glyn Morrill	61
Jean Carletta	110	Nicolas Nicolov	217
John Carroll	217	Gertjan van Noord	126
John Chen	188	Franz Josef Och	71
Dimitris N. Christodoulakis	134	Ian M. O'Neill	23
Fabio Ciravegna	102	Giorgos S. Orphanos	134
Eric de la Clergerie	150	Richard Power	9
Myriam Eckert	37	Erik F. Tjong Kim Sang	173
Nikos Fakotakis	158	Satoshi Sekine	196
Tim Fernando	86	Olga Shaumyan	217
Therese Firmin	77	Hiroyuki Shinnou	180
Fumiyo Fukumoto	209	Martine Smets	217
Dale Gerdemann	126	Efstathios Stamatatos	158
Claire Grover	1	Suzanne Stevenson	45
Lynette Hirschman	77	Michael Strube	37
David House	77	Beth Sundheim	77
Hitoshi Isahara	196	Yoshimi Suzuki	209
Christian Jacquemin	15	Hozumi Tanaka	94
Frank Keller	30	Simone Teufel	110
Mathias Kirsten	142	Takenobu Tokunaga	94
Gary Klein	77	Kiyotaka Uchimoto	196
George Kokkinakis	158	Jorn Veenstra	173
Maria Lapata	30	K. Vijay-Shanker	188
Alberto Lavelli	102	Manuel Vilares	150
Rila Mandala	94	Atro Voutilainen	204
Inderjeet Mani	77	David Weir	217
Scott McDonald	30	Maria Wolters	142

Student Author Index

Donna K. Byron	229	Daniel S. Paiva	237
Patrick Caudal	233	Justin Picard	241
Gabriela Cavaglià	225	Veit Reuer	253
Sophie Johansson Kokkinakis	245	Joel R. Tetreault	229
Dimitrios Kokkinakis	245	Aline Villavicencio	261
Patrice Lopez	249	Yan Zuo	257

Poster Author Index

Itziar Aduriz	265	Katsutoshi Ibushi	271
Izaskun Aldezabal	265	Hisao Imai	271
Inaki Alegria	265	Adam Kilgarriff	277
Olatz Ansa	265	Peter Kleiweg	281
Hans Aregenton	267	Torbjörn Lager	279
Jose Mari Arriola	265	Guido Minnen	269
Yvonne Canning	269	Kaili Müürisep	291
John Carroll	269	John Nerbonne	281
Marcela Charfuelan	287	Chikashi Nobata	271
Nigel Collier	271	Norihiro Ogata	271
Alexander Da Costa	265	Tomoko Ohta	271
Crit Cremers	273	Hyun Seok Park	271
Siobhan Devlin	269	Darren Pearce	269
Nerea Ezeiza	265	Frank Schilder	283
Anke Feldhaus	267	Tateshi Sekimizu	271
Johann Gamper	275	Mark Stevenson	285
Maria G. Gancedo	287	John Tait	269
José Relano Gil	287	Daniel Tapias	287
Wilbert Heeringa	281	Yuka Tateishi	271
Luis A. Hernández	287	Jun-ichi Tsujii	271

EACL '99 Programme

June 9, 1999

- 0830 - 0930 Registration
0930 - 0940 Welcome
0940 - 1020 Main session 1
Auditorium 2
Andrei Mikheev, Marc Moens and Claire Grover: Named Entity Recognition without Gazetteers
Auditorium 3
Richard Power: Generating Referring Expressions with a Unification Grammar
- 1020 - 1100 Main session 2
Auditorium 2
Didier Bourigault and Christian Jacquemin: Term Extraction and Term Clustering: An Integrated Platform for Computer-Aided Terminology
Auditorium 3
Ian M. O'Neill and Michael McTear: An Object-Oriented Approach to the Design of Dialogue Management Functionality
- 1100 - 1130 Break
- 1130 - 1230 Invited talk
Auditorium 2
Bruce Croft: Language Models for Information Retrieval
- 1230 - 1400 Lunch
- 1400 - 1430 Student session 1
Auditorium 2
Gabriela Cavaglià: The Development of Lexical Resources for Information Extraction from Text Combining WordNet and Dewey Decimal Classification
Auditorium 3
Donna K. Byron and Joel R. Tetreault: A Flexible Architecture for Reference Resolution
- 1430 - 1510 Main session 3
Auditorium 2
Maria Lapata, Scott McDonald and Frank Keller: Determinants of Adjective-Noun Plausibility
Auditorium 3
Miriam Eckert and Michael Strube: Resolving Discourse Deictic Anaphora in Dialogues
- 1510 - 1550 Break
- 1550 - 1620 Student session 2
Auditorium 2
Patrick Caudal: Result States and the Lexicon: the Proper Treatment of Event Structure
Auditorium 3
Daniel S. Paiva: Investigating NLG Architectures: Taking Style into Consideration
- 1620 - 1700 Main session 4
Auditorium 2
Susanne Stevenson and Paola Merlo: Automatic Verb Classification Using Distributions of Grammatical Features
Auditorium 3
Pierre Boullier: Chinese Numbers, MIX, Scrambling and Range Concatentation Grammars

June 10, 1999

- 0900 - 0940 Main session 1
Auditorium 2
Glyn Morrill: Geometry of Lexico-Syntactic Interaction
Auditorium 3
Franz Josef Och: An Efficient Method for Determining Bilingual Word Classes
- 0940 - 1020 Main session 2
Auditorium 2
Inderjeet Mani, Therese Firmin, David House, Gary Klein, Beth Sundheim, Lynette Hirschman: The TIPSTER SUMMAC Text Summarization Evaluation
Auditorium 3
Tim Fernando: Ambiguous Propositions Typed
- 1020 - 1050 Break
- 1050 - 1150 Invited talk
Auditorium 2
Wolfgang Wahlster: Deep Processing of Shallow Structures: The Robust Integration of Speech, Language and Translation Technology for Intelligent Interface Agents
- 1150 - 1230 Main session 3
Auditorium 2
Rila Mandala, Takenobu Tokunaga and Hozumi Tanaka: Complementing WordNet with Roget's and Corpus-based Thesauri for Information Retrieval
Auditorium 3
Fabio Ciravegna and Alberto Lavelli: Full Text Parsing Using Cascades of Rules: An Information Extraction Perspective
- 1230 - 1400 Lunch
- 1400 - 1430 Student session 1
Auditorium 2
Justin Picard: Finding Content-bearing Terms using Term Similarities
Auditorium 3
Dimitrios Kokkinakis and Sophie Johansson Kokkinakis: A Cascaded Finite-State Parser for Syntactic Analysis of Swedish
- 1430 - 1510 Main session 4
Auditorium 2
Simone Teufel, Jean Carletta and Marc Moens: An Annotation Scheme for Discourse Level Argumentation in Research Articles
Auditorium 3
Thorsten Brants: Cascaded Markov Models
- 1510 - 1550 Break
- 1550 - 1620 Student session 2
Auditorium 2
Patrice Lopez: Repair Strategies for Lexicalized Tree Grammars
Auditorium 3
Veit Reuer: Dialogue Processing in a CALL-System
- 1620 - 1700 Main session 5
Auditorium 2
Dale Gerdemann and Gertjan van Noord: Transducers from Rewrite Rules with Backreferences
Auditorium 3
Giorgos Orphanos and Dimoitis Christodoulakis: POS Disambiguation and Unknown Word Guessing with Decision Trees

June 11, 1999

- 0900 - 0940 Main session 1
Auditorium 2
Maria Wolters and Mathias Kirsten: Exploring the Use of Linguistic Features in Domain and Genre Classification
Auditorium 3
Miguel Alonso, David Cabrero, Eric de la Clergerie and Manuel Vilares: Tabular Algorithms for TAG Parsing
- 0940 - 1020 Main session 2
Auditorium 2
Efsthathios Stamatatos, Nikos Fakotakis and George Kokkinakis: Automatic Authorship Attribution
Auditorium 3
Guido Minnen: Selective Magic HPSG Parsing
- 1020 - 1050 Break
- 1050 - 1150 Poster and Demo session
- 1150 - 1230 Main session 3
Auditorium 2
Erik Tjong Kim Sang and Jorn Veenstra: Representing Text Chunks
Auditorium 3
Hiroyuki Shinnou: Detection of Japanese Homophone Errors by a Decision List Including a Written Word as a Default Evidence
- 1230 - 1400 Lunch
- 1400 - 1440 Main session 4
Auditorium 2
John Chen, Srinivas Bangalore and K. Vijay-Shanker: New Models for Improving Supertag Disambiguation
Auditorium 3
Kiyotaka Uchimoto, Satoshi Sekine and Hitoshi Isahara: Japanese Dependency Structure Analysis based on Maximum Entropy Models
- 1440 - 1520 Main session 5
Auditorium 2
Atro Voutilainen: An Experiment on the Upper Bound of Interjudge Agreement: The Case of Tagging
Auditorium 3
Fumiyo Fukumoto and Yoshimi Suzuki: Word Sense Disambiguation in Untagged Text Based on Term Weight Learning
- 1520 - 1535 Break
- 1535 - 1615 Business meeting (incl. TEI)
Auditorium 2

Posters

Izaskun Aldezabal, Inaki Alegria, Olatz Ansa, Jose Mari Arriola, Nerea Ezeiza, Itziar Aduriz and Alexander Da Costa: Designing Spelling Correctors for Inflected Languages Using Lexical Transducers

Hans Argenton and Anke Feldhaus: The Treegram Index—An Efficient Technique for Retrieval in Linguistic Treebanks

John Carroll, Guido Minnen, Darren Pearce, Yvonne Canning, Siobhan Devlin and John Tait: Simplifying Text for Language-Impaired Readers

Nigel Collier, Hyun Seok Park, Norihiro Ogata, Yuka Tateishi, Chikashi Nobata, Tomoko Ohta, Tateshi Sekimizu, Hisao Imai, Katsutoshi Ibushi, Jun-ichi Tsujii: The GENIA Project: Corpus-Based Knowledge Acquisition and Information Extraction from Genome Research Papers

Crit Cremers: A Note on Categorical Grammar, Disharmony and Permutation

Johann Gamper: Encoding a Parallel Corpus for Automatic Terminology Extraction

José Relano Gil, Daniel Tapias, Maria C. Gancedo, Marcela Charfuelan and Luis A. Hernández: Robust and Flexible Mixed-Initiative Dialogue for Telephone Services

Adam Kilgarriff: 95% Replicability for Manual Word Sense Tagging

Torbjörn Lager: μ -TBL: A Small, Extendible Transformation-Based Learner

Kaili Müürisep: Determination of Syntactic Functions in Estonian Constraint Grammar

John Nerbonne, Wilbert Heeringa and Peter Kieweg: Comparison and Classification of Dialects

Frank Schilder: Pointing to Events

Mark Stevenson: A Corpus-Based Method for Deriving Lexical Mappings

Tutorials

The tutorials take place on June 8th, 1999.

Practical Text Mining

0930–1300

Lecturer: **Ronen Feldman** (Bar-Ilan University)

The information age has made it easy to store large amounts of data. The proliferation of documents available on the Web, on corporate intranets, on news wires, and elsewhere is overwhelming. However, while the amount of data available to us is constantly increasing, our ability to absorb and process this information remains constant. Search engines only exacerbate the problem by making more and more documents available in a matter of a few key strokes. Text Mining is a new and exciting research area that tries to solve the information overload problem by using techniques from data mining, machine learning, NLP, IR and knowledge management. Text Mining involves the preprocessing of document collections (text categorization, term extraction), the storage of the intermediate representations, the techniques to analyze these intermediate representations (distribution analysis, clustering, trend analysis, association rules etc) and visualization of the results. In this tutorial we will present the general theory of Text Mining and will demonstrate several systems that use these principles to enable interactive exploration of large textual collections. We will present a general architecture for text mining and will outline the algorithms and data structures behind the systems. Special emphasis will be given to efficient algorithms for very large document collections, tools for visualizing such document collections, the use of intelligent agents to perform text mining on the internet, and the use information extraction to better capture the major themes of the documents. The Tutorial will cover the state of the art in this rapidly growing area of research. Several real world applications of text mining will be presented.

Natural Language Learning with the Maximum Entropy Framework

1400–1730

Lecturer: **Adwait Ratnaparkhi** (IBM TJ Watson Research Center)

“Corpus-based” approaches to natural language processing (NLP), also known as “statistical” or “machine learning” approaches, have become popular in recent years due to the availability of large, annotated corpora. This tutorial will discuss how to implement a corpus-based NLP tool with the maximum entropy framework. We will first describe the maximum entropy framework, and then discuss its application to several problems, including sentence boundary detection, part-of-speech tagging, prepositional phrase attachment, parsing, and text categorization. Our experience has shown that this framework yields consistently high accuracies, requires relatively “knowledge-poor” informants, and is highly reusable across tasks. A general outline of the tutorial:

1. What is the maximum entropy framework ?
2. Using it for integrating diverse sources of evidence:
 - Sentence Boundary Detection
 - Part of Speech tagging
 - Parsing
3. Using it with and without annotated data:
 - Prepositional Phrase Attachment
4. Comparing its performance with other learning techniques:
 - Text Categorization
 - Incremental vs. Frequency-based feature selection
 - Decision Trees (C5.0)
5. Summary: Advantages and Disadvantages of the framework

Natural language generation (NLG) systems produce understandable texts in English and other human languages from some underlying non-linguistic representation of information. NLG systems combine knowledge about language and the application domain to automatically produce documents, reports, explanations, help messages, and other kinds of texts.

In this tutorial we describe NLG from an applied system-building perspective; that is, we will explain how NLG systems are built. Our presentation will be based on a popular architectural model which encompasses the three stages of text planning, microplanning, and realisation. We will also give examples of current NLG applications; discuss when NLG technology is and is not appropriate; and explore how NLG can be integrated into multimedia, hypertext, and speech systems.

This tutorial should be useful for managers, implementors, and researchers. For managers, it will provide a broad overview of the field and what is possible today; for implementors, it will provide a realistic assessment of available techniques; and for researchers, it will highlight the issues that are important in current applied NLG projects.

NLP makes extensive use of dictionaries, but frequently proceeds in ignorance of lexicography. Even on computational lexicography projects, workers rarely have a background in lexicography. If workers in NLP were better informed about how dictionaries were produced, the purposes they were designed to serve, and what distinguished good dictionaries from bad ones, they would be better placed to choose a dictionary and to exploit the information it contained.

The tutorial will describe the goals and practices of corpus-based dictionary-making. Participants will gain an appreciation of the kinds of judgements that lexicographers need to make every day and the criteria they use. The tutorial will include practical exercises.

We shall also discuss models for lexicography/NLP collaboration, including SENSEVAL, the recent evaluation exercise for word sense disambiguation programs.